

JOURNAL OF COMPUTER AND SYSTEM SCIENCES 26, 171-196 (1983)

On the Equality of Grammatical Families

SEYMOUR GINSBURG*

*Computer Science Department, University of Southern California,
Los Angeles, California 90089*

JONATHAN GOLDSTINE†

*Computer Science Department, The Pennsylvania State University,
University Park, Pennsylvania 16802*

AND

EDWIN H. SPANIER

*Mathematics Department, University of California,
Berkeley, California 94720*

Received March 20, 1982; revised August 1982

A decision procedure involving the testing for containment is presented for determining whether two grammatical families (families generated by context-free grammar forms) are equal. It is also shown that every nontrivial grammatical family which is a proper subset of the family of context-free languages can be constructed in a unique canonical way from the family of regular sets by certain operations. Thus, two such grammatical families are equal iff their canonical representations are identical.

INTRODUCTION

In 1975, Cremers and Ginsburg [3] introduced the concept of a family (called a *grammatical family*) of languages generated by context-free grammars which resemble a fixed master grammar. Although many papers on grammatical families have appeared since then (see [6, 8] for bibliographies), some basic problems have remained open for some time, including the problem of determining whether two

* This author was supported in part by the National Science Foundation under Grants MCS 77-22323 and MCS 7925004.

† This author was supported in part by the National Science Foundation under Grant MCS 76-10076A01.

grammatical families are equal. This suggests that additional tools are needed for the analysis of grammatical families. In an earlier paper [7], the authors developed such tools by taking an algebraic approach to the analysis of grammatical families, and establishing the

PRIME DECOMPOSITION THEOREM. *Every grammatical family can be expressed as a minimal sum of products of prime grammatical families in a unique way.* (The relevant terms will be defined in the next section)

These tools reduce the problem of determining whether one grammatical family includes another to the special case in which the first family is prime and the second is a product of primes. This paper resolves that special case to obtain a general decision procedure for the inclusion of one grammatical family in another, and hence also for the equality of grammatical families.¹ It is also proven here that each nontrivial grammatical family properly contained in the context-free family has a unique canonical expression in terms of the family \mathcal{R} of regular languages and certain operations on families. Thus, two grammatical families are equal if and only if their canonical representations are identical. Since the canonical expression for a grammatical family refines its minimal prime decomposition, this result extends the prime decomposition theorem to give a complete analysis of a grammatical family into unique components.

The paper is divided into three sections and an Appendix. Section 1 contains definitions, notation, and some preliminary lemmas. Section 2 contains the decision procedures for containment and equality of grammatical families. Section 3 establishes canonical forms for these families. And the Appendix contains the proof of one of the lemmas from Section 1.

1. PRELIMINARIES

In this section, we first review some material about grammatical families. The reader is referred to [7] for further details. We then introduce notation to describe certain types of these families. Finally, we establish a number of technical lemmas and corollaries needed for the development (in Sections 2 and 3) of our major results.

Let V_∞ be a fixed infinite universe of symbols and Σ_∞ a subset of V_∞ such that Σ_∞ and $V_\infty - \Sigma_\infty$ are both infinite. All nonterminals used in grammars² are to be elements of $V_\infty - \Sigma_\infty$ and all terminals to be elements of Σ_∞ . Thus, $G = (V, \Sigma, P, \sigma)$ is a grammar if Σ is a finite subset of Σ_∞ , V is a finite subset of V_∞ , $(V - \Sigma) \subseteq (V_\infty - \Sigma_\infty)$, σ is in $V - \Sigma$, and P is a finite subset of $(V - \Sigma) \times V^*$.

We shall view a grammar in two different senses. The first is the customary one of

¹ In 1978, without spelling out the details, we announced the solution to the equality problem [6]. More recently, Meera Blattner also announced a solution to this problem [1], but by a different method, one not based on a systematic decomposition theory for grammatical families.

² By a *grammar* is always meant a context-free grammar.

a language-generating device. The second is as a master grammar defining a family of grammars, each "looking like" the master one. (The latter is referred to in the literature as the "grammar form" sense.) A grammar G determines grammars which "look like" G by the notion of

DEFINITION. An *interpretation* of a grammar $G = (V, \Sigma, P, \sigma)$ is a 5-tuple $I = (\mu_I, V_I, \Sigma_I, P_I, S_I)$, where $G_I = (V_I, \Sigma_I, P_I, S_I)$ is a grammar and μ_I is a substitution on V^* satisfying the following conditions:

- (i) $\mu_I(\xi) \subseteq V_I - \Sigma_I$ for all ξ in $V - \Sigma$,
- (ii) $\mu_I(a)$ is a finite subset of Σ_I^* for all a in Σ ,
- (iii) $\mu_I(\xi) \cap \mu_I(\eta) = \emptyset$ for all $\xi \neq \eta$ in $V - \Sigma$,
- (iv) S_I is in $\mu_I(\sigma)$, and
- (v) $P_I \subseteq \mu_I(P)$, where $\mu_I(\xi \rightarrow w) = \{\alpha \rightarrow y \mid \alpha \text{ in } \mu_I(\xi), y \text{ in } \mu_I(w)\}$ and $\mu_I(P) = \bigcup_{p \text{ in } P} \mu_I(p)$.

Intuitively, G_I is supposed to "look like" G .

The central object of our study is the family of languages arising from the different interpretation grammars.

DEFINITION. A collection \mathcal{L} of languages is a *grammatical family* if $\mathcal{L} = \{L(G_I) \mid I \text{ an interpretation of } G\}$ for some grammar G , called its *grammar form*. \mathcal{L} is *trivial* if it contains only finite languages, and *nontrivial*, otherwise. \mathcal{L} is *proper* if it is nontrivial and is properly contained in the family \mathcal{L}_{CF} of all context-free languages.

The only trivial grammatical families are $\mathcal{L}_\emptyset = \{\emptyset\}$, $\mathcal{L}_\varepsilon = \{\emptyset, \{\varepsilon\}\}$,³ and $\mathcal{L}_{fin} = \{L \mid L \text{ finite}\}$. There are decision procedures for determining whether a grammatical family (given by a grammar form) equals \mathcal{L}_\emptyset , or \mathcal{L}_ε , or \mathcal{L}_{fin} , or \mathcal{L}_{CF} [4]. The nontrivial grammatical families are all full principal semi-AFLs⁴ [3].

DEFINITION. The *sum* of two nonempty families of languages \mathcal{L} and \mathcal{L}' , denoted by $\mathcal{L} \oplus \mathcal{L}'$, is $\{L \cup L' \mid L \text{ in } \mathcal{L}, L' \text{ in } \mathcal{L}'\}$; and their *product*, denoted by $\mathcal{L} \odot \mathcal{L}'$, is $\{\bigcup_{i=1}^n L_i L'_i \mid n \geq 1, \text{ each } L_i \text{ in } \mathcal{L} \text{ and } L'_i \text{ in } \mathcal{L}'\}$.

Both sum and product are obviously associative, so parentheses may be omitted. Furthermore,

$$(\mathcal{L}_1 \oplus \mathcal{L}_2) \odot \mathcal{L}_3 = (\mathcal{L}_1 \odot \mathcal{L}_3) \oplus (\mathcal{L}_2 \odot \mathcal{L}_3)$$

³ ε denotes the empty word.

⁴ A *full semi-AFL* is a family of languages containing a nonempty language and closed under homomorphism, inverse homomorphism, intersection with regular sets, and union. A *full AFL* is a full semi-AFL which is closed under product and $*$. If \mathcal{L} is a family of languages, let $\mathcal{S}(\mathcal{L})$ and $\mathcal{F}(\mathcal{L})$ be the smallest full semi-AFL and smallest full AFL containing \mathcal{L} . A full semi-AFL \mathcal{L} is called *full principal* if there exists a language L such that $\mathcal{L} = \mathcal{S}(\{L\})$, usually written $\mathcal{L} = \mathcal{S}(L)$. For further details, the reader is referred to [5].

and

$$\mathcal{L}_3 \odot (\mathcal{L}_1 \oplus \mathcal{L}_2) = (\mathcal{L}_3 \odot \mathcal{L}_1) \oplus (\mathcal{L}_3 \odot \mathcal{L}_2),$$

provided that \mathcal{L}_1 and \mathcal{L}_2 both contain the empty language. Thus, product distributes over sum for grammatical families.

DEFINITION. An expression of the form

$$(\mathcal{L}_{11} \odot \cdots \odot \mathcal{L}_{1n_1}) \oplus \cdots \oplus (\mathcal{L}_{m1} \odot \cdots \odot \mathcal{L}_{mn_m}),$$

$m \geq 1$, each $n_i \geq 1$, is called a *sum-of-products* expression, with each $\mathcal{L}_{i1} \odot \cdots \odot \mathcal{L}_{in_i}$ a *summand*. The expression is *minimal* if, for every i and j , the deletion of \mathcal{L}_{ij} produces an expression which designates a strictly smaller family.

Since summation is commutative, the ordering of the summands in an expression will be ignored. In other words, we shall consider two expressions to be identical if they differ only by a permutation of their summands.

DEFINITION. A grammatical family \mathcal{L} is *prime* (respectively, *additively prime*) if, for every pair of grammatical families \mathcal{L}_1 and \mathcal{L}_2 such that $\mathcal{L} \subseteq \mathcal{L}_1 \odot \mathcal{L}_2$ (respectively, $\mathcal{L} \subseteq \mathcal{L}_1 \oplus \mathcal{L}_2$), either $\mathcal{L} \subseteq \mathcal{L}_1$ or $\mathcal{L} \subseteq \mathcal{L}_2$.

If a family is prime, then it is additively prime. In fact, a family is additively prime if and only if it is the finite product of one or more prime families [7].

The following theorem is proved in [7]:

PRIME DECOMPOSITION THEOREM. Every grammatical family is a unique minimal sum of products of prime grammatical families.

The collection of proper grammatical families is the smallest collection containing the family \mathcal{R} of regular sets and closed under the operations \oplus , \odot , \mathcal{F} , and a ternary operator \mathcal{E} defined as follows [7]:

DEFINITION. A grammar $G = (V, \Sigma, P, \sigma)$ is a *split linear* grammar if the right-hand side of every production in P is in $A(V - \Sigma) \cup C \cup (V - \Sigma)B$ for some disjoint subsets A , B , and C of Σ . In such a case, we use the notation $G = (V, A \cup C \cup B, P, \sigma)$. For all families \mathcal{L}_1 , \mathcal{L}_2 , \mathcal{L}_3 of languages, $\mathcal{E}(\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3)$ is the family of all languages of the form $\tau(L)$, where $L = L(G)$ for some split linear grammar $G = (V, A \cup C \cup B, P, \sigma)$ and τ is a substitution on $(A \cup C \cup B)^*$ such that $\tau(x)$ is in \mathcal{L}_1 , \mathcal{L}_2 , or \mathcal{L}_3 if x is in A , C , or B , respectively.

Clearly, $\mathcal{E}(\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3) \subseteq \mathcal{E}(\mathcal{L}'_1, \mathcal{L}'_2, \mathcal{L}'_3)$ if $\mathcal{L}_i \subseteq \mathcal{L}'_i$ for each i .

Notation. The following five symbols, with or without subscripts or primes, will have the indicated meanings:

\mathcal{G} : a proper grammatical family, i.e., a grammatical family such that $\mathcal{R} \subseteq \mathcal{G} \subseteq \mathcal{L}_{CF}$.

- \mathcal{A} : an additively prime proper grammatical family.
- \mathcal{P} : a prime proper grammatical family.
- \mathcal{F} : a proper grammatical family of the form $\mathcal{F}(\mathcal{G})$.
- \mathcal{E} : a proper grammatical family of the form $\mathcal{E}(\mathcal{G}, \mathcal{A}, \mathcal{G}')$.

It is shown in [7] that a proper grammatical family is prime if and only if it has the form $\mathcal{F}(\mathcal{G})$ or $\mathcal{E}(\mathcal{G}, \mathcal{A}, \mathcal{G}')$. It is shown later in this section (Lemma 1.10) that these types do not overlap, i.e., for no \mathcal{F} and \mathcal{E} is it true that $\mathcal{F} = \mathcal{E}$.

Note that $\mathcal{R} \odot \mathcal{G} = \mathcal{G} \odot \mathcal{R} = \mathcal{G}$ for all \mathcal{G} , i.e., the family \mathcal{R} of regular sets serves as a multiplicative identity for the collection of all proper grammatical families. Consequently, the convention will be adopted that $\mathcal{G}_i \odot \mathcal{G}_{i+1} \odot \cdots \odot \mathcal{G}_j$ denotes \mathcal{R} when $j < i$.

We shall write $(\mathcal{G}'_1, \mathcal{G}'_2) \subseteq (\mathcal{G}_1, \mathcal{G}_2)$ as an abbreviation for $\mathcal{G}'_1 \subseteq \mathcal{G}_1$ and $\mathcal{G}'_2 \subseteq \mathcal{G}_2$, and shall say that $(\mathcal{G}'_1, \mathcal{G}'_2)$ and $(\mathcal{G}_1, \mathcal{G}_2)$ are *incomparable* if neither $(\mathcal{G}'_1, \mathcal{G}'_2) \subseteq (\mathcal{G}_1, \mathcal{G}_2)$ nor $(\mathcal{G}_1, \mathcal{G}_2) \subseteq (\mathcal{G}'_1, \mathcal{G}'_2)$ is true.

The rest of the section is concerned with a number of technical results about the \mathcal{E} operator. We start with some elementary facts (four lemmas and three corollaries).

For all $\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3$, $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$ is the family of all finite unions of languages of the form

$$(*) \quad \bigcup_{w \text{ in } R} \tau_1(w) L_2 \tau_3(w^R),$$

where R is in \mathcal{R} , L_2 is in \mathcal{G}_2 , and τ_i is a \mathcal{G}_i -substitution.⁵

Lemmas 1.1 and 1.2 are derived from (*). Details are in [7].

1.1 LEMMA. $\mathcal{E}(\mathcal{G}, \mathcal{G}_1 \oplus \cdots \oplus \mathcal{G}_n, \mathcal{G}') = \mathcal{E}(\mathcal{G}, \mathcal{G}_1, \mathcal{G}') \oplus \cdots \oplus \mathcal{E}(\mathcal{G}, \mathcal{G}_n, \mathcal{G}')$.

1.2 LEMMA. $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) = \mathcal{E}(\mathcal{F}(\mathcal{G}_1), \mathcal{G}_2, \mathcal{F}(\mathcal{G}_3))$.

Lemmas 1.3 and 1.4 are also derived from (*).

1.3 LEMMA. $\mathcal{G}_1 \odot \mathcal{G}_2 \odot \mathcal{G}_3 \subseteq \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$.

Proof. This follows from (*) by letting $R = \{a\}$, where a is a single symbol. ■

1.4 LEMMA. $\mathcal{E}(\mathcal{G}_1, \mathcal{E}(\mathcal{G}'_1, \mathcal{G}_2, \mathcal{G}'_3), \mathcal{G}_3) \subseteq \mathcal{E}(\mathcal{G}_1 \oplus \mathcal{G}'_1, \mathcal{G}_2, \mathcal{G}_3 \oplus \mathcal{G}'_3)$.

Proof. By (*), each language in the left-hand family is a finite union of languages of the form

$$L = \bigcup_{u \text{ in } R} \bigcup_{v \text{ in } R'} \tau_1(u) \tau'_1(v) L_2 \tau'_3(v^R) \tau_3(u^R),$$

where R, R' are in \mathcal{R} , L_2 is in \mathcal{G}_2 , and $\tau_i(\tau'_i)$ is a \mathcal{G}_i (\mathcal{G}'_i -) substitution. Without loss

⁵ A substitution τ is a \mathcal{G} -substitution if $\tau(x)$ is in \mathcal{G} for each symbol x .

of generality we may assume that R and R' have disjoint alphabets, so τ_i and τ'_i can be extended to a single $\mathcal{E}_i \oplus \mathcal{E}'_i$ -substitution σ_i on the union of their alphabets. Then

$$L = \bigcup_{w \text{ in } RR'} \sigma_1(w) L_2 \sigma_3(w^R).$$

By (*), finite unions of languages of this form are in $\mathcal{E}(\mathcal{E}_1 \oplus \mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}_3 \oplus \mathcal{E}'_3)$. ■

1.5 COROLLARY. $\mathcal{E}_1 \odot \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3) \odot \mathcal{E}_3 = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3)$.

Proof.

$$\begin{aligned} & \mathcal{E}_1 \odot \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3) \odot \mathcal{E}_3 \\ & \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3), \mathcal{E}_3) \quad \text{by Lemma 1.3,} \\ & \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3) \quad \text{by Lemma 1.4 (since } \mathcal{E}_i \oplus \mathcal{E}_i = \mathcal{E}_i). \end{aligned}$$

The reverse inclusion is immediate. ■

1.6 COROLLARY. If $(\mathcal{E}'_1, \mathcal{E}'_3) \subseteq (\mathcal{E}_1, \mathcal{E}_3)$, then

$$\mathcal{E}(\mathcal{E}'_1, \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3), \mathcal{E}'_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}(\mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}'_3), \mathcal{E}_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3).$$

Proof. Suppose $(\mathcal{E}'_1, \mathcal{E}'_3) \subseteq (\mathcal{E}_1, \mathcal{E}_3)$. Then $\mathcal{E}_1 \oplus \mathcal{E}'_1 = \mathcal{E}_1$ and $\mathcal{E}_3 \oplus \mathcal{E}'_3 = \mathcal{E}_3$. Hence, by Lemma 1.3,

$$\mathcal{E}(\mathcal{E}'_1, \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3), \mathcal{E}'_3) \subseteq \mathcal{E}(\mathcal{E}_1 \oplus \mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}_3 \oplus \mathcal{E}'_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3)$$

and

$$\mathcal{E}(\mathcal{E}_1, \mathcal{E}(\mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}'_3), \mathcal{E}_3) \subseteq \mathcal{E}(\mathcal{E}_1 \oplus \mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}_3 \oplus \mathcal{E}'_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3).$$

The reverse inclusions are immediate. ■

1.7 COROLLARY. If $(\mathcal{E}'_1, \mathcal{E}'_3) \subseteq (\mathcal{E}_1, \mathcal{E}_3)$, then

$$\mathcal{E}(\mathcal{E}_1, \mathcal{E}'_1 \odot \mathcal{E}_2, \mathcal{E}_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2 \odot \mathcal{E}'_3, \mathcal{E}_3) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3).$$

Proof. Suppose $(\mathcal{E}'_1, \mathcal{E}'_3) \subseteq (\mathcal{E}_1, \mathcal{E}_3)$. By symmetry, it suffices to show the equality of $\mathcal{E}(\mathcal{E}_1, \mathcal{E}'_1 \odot \mathcal{E}_2, \mathcal{E}_3)$ and $\mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3)$. Now

$$\begin{aligned} \mathcal{E}(\mathcal{E}_1, \mathcal{E}'_1 \odot \mathcal{E}_2, \mathcal{E}_3) & \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{E}'_1 \odot \mathcal{E}_2 \odot \mathcal{E}'_3, \mathcal{E}_3) \\ & \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{E}(\mathcal{E}'_1, \mathcal{E}_2, \mathcal{E}'_3), \mathcal{E}_3) \quad \text{by Lemma 1.3,} \\ & \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{E}_2, \mathcal{E}_3) \quad \text{by Corollary 1.6.} \end{aligned}$$

The reverse inclusion is immediate. ■

The remaining results are somewhat deeper, and are needed for the containment decision procedure in the next section.

1.8 LEMMA. $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \subseteq \mathcal{F}(\mathcal{G})$ iff $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \subseteq \mathcal{G}$.

Proof. Since $\mathcal{G} \subseteq \mathcal{F}(\mathcal{G})$, the if direction is immediate. To show the converse, suppose that $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \subseteq \mathcal{F}(\mathcal{G})$. Let L be in $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$ and

$$L' = \{a^n w b^n \mid n \geq 1, w \text{ in } L\},$$

where a and b are new letters. Then L' is in

$$\begin{aligned} & \mathcal{E}(\mathcal{R}, \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3), \mathcal{R}) \\ & \subseteq \mathcal{E}(\mathcal{G}_1, \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3), \mathcal{G}_3) \quad \text{since } \mathcal{R} \subseteq \mathcal{G}_i \text{ for each } i, \\ & = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \quad \text{by Corollary 1.5,} \\ & \subseteq \mathcal{F}(\mathcal{G}). \end{aligned}$$

But by [2, Theorem 1], given a full semi-AFL \mathcal{L} , a language of the form L' is in $\mathcal{F}(\mathcal{L})$ iff it is in \mathcal{L} . Hence, L' is in \mathcal{G} , so L is in \mathcal{G} . Therefore $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \subseteq \mathcal{G}$. ■

1.9 LEMMA. If $\mathcal{A}_1 \odot \mathcal{A}_2 \subseteq \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, then either $\mathcal{A}_1 \subseteq \mathcal{F}_1$, $\mathcal{A}_2 \subseteq \mathcal{F}_2$, or $\mathcal{A}_1 \odot \mathcal{A}_2 \subseteq \mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$.

Proof. The argument is involved, and is relegated to the Appendix. ■

1.10 LEMMA. For all \mathcal{F} and \mathcal{E} , $\mathcal{F} \neq \mathcal{E}$.

Proof. Suppose $\mathcal{F} = \mathcal{E}$ for some \mathcal{F} and \mathcal{E} . Since \mathcal{E} is prime, it follows from [7, Corollary 4 of the prime decomposition theorem] that \mathcal{E} has one of the following forms:

- (i) $\mathcal{E} = \mathcal{R}$,
- (ii) $\mathcal{E} = \mathcal{F}(\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n)$, with $n \geq 1$ and each $\mathcal{P}_i \neq \mathcal{E}$, or
- (iii) $\mathcal{E} = \mathcal{E}(\mathcal{P}_1, \mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}, \mathcal{P}_n)$, with $n \geq 3$ and each $\mathcal{P}_i \neq \mathcal{E}$.

The smallest value that \mathcal{E} can be is $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})$, which is the family of linear languages. Therefore (i) cannot occur. Thus (ii) or (iii) holds. In either case, $\mathcal{P}_i \subseteq \mathcal{E}$ for each i . Thus, $\mathcal{E} \not\subseteq \mathcal{P}_i$ for each i . Assume (ii) holds. Then $\mathcal{E} \subseteq \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$ by Lemma 1.8. Since \mathcal{E} is prime and $\mathcal{E} \not\subseteq \mathcal{P}_i$ for each i , this is a contradiction. Thus (iii) holds. Then

$$\begin{aligned} \mathcal{E} \odot \mathcal{E} &= \mathcal{F} \odot \mathcal{F} = \mathcal{F} = \mathcal{E} \\ &= \mathcal{E}(\mathcal{P}_1, \mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}, \mathcal{P}_n) \\ &\subseteq \mathcal{E}(\mathcal{F}(\mathcal{P}_1), \mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}, \mathcal{F}(\mathcal{P}_n)). \end{aligned}$$

Since $\mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}$ is the product of one or more prime families, it is additively prime. By Lemma 1.9, either $\mathcal{E} \subseteq \mathcal{F}(\mathcal{P}_1)$, $\mathcal{E} = \mathcal{E} \odot \mathcal{E} \subseteq \mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}$, or $\mathcal{E} \subseteq \mathcal{F}(\mathcal{P}_3)$. Since $\mathcal{E} \not\subseteq \mathcal{P}_1$ and $\mathcal{E} \not\subseteq \mathcal{P}_3$, $\mathcal{E} \not\subseteq \mathcal{F}(\mathcal{P}_1)$ and $\mathcal{E} \not\subseteq \mathcal{F}(\mathcal{P}_3)$ by Lemma 1.8. Hence $\mathcal{E} \subseteq \mathcal{P}_2 \odot \cdots \odot \mathcal{P}_{n-1}$. As in (ii), this leads to a contradiction. ■

Although not needed for our development, we note in passing that Lemma 1.10 can be extended to show that $\mathcal{F}(\mathcal{E}) \neq \mathcal{E}(\mathcal{E}_1, \mathcal{E}_3, \mathcal{E}_2)$ for all $\mathcal{E}, \mathcal{E}_1, \mathcal{E}_2$, and \mathcal{E}_3 . Indeed, suppose $\mathcal{F}(\mathcal{E}) = \mathcal{E}(\mathcal{E}_1, \mathcal{E}_3, \mathcal{E}_2)$ for some $\mathcal{E}, \mathcal{E}_1, \mathcal{E}_2$, and \mathcal{E}_3 . By the prime decomposition theorem and by [7, Proposition C], \mathcal{E}_3 has the form $\mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_m$, $m \geq 1$. Then

$$\begin{aligned}\mathcal{F}(\mathcal{E}) &= \mathcal{E}(\mathcal{E}_1, \mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_m, \mathcal{E}_2) \\ &= \mathcal{E}(\mathcal{E}_1, \mathcal{A}_1, \mathcal{E}_2) \oplus \cdots \oplus \mathcal{E}(\mathcal{E}_1, \mathcal{A}_m, \mathcal{E}_2) \quad \text{by Lemma 1.1.}\end{aligned}$$

Since $\mathcal{F}(\mathcal{E})$ is prime, $\mathcal{F}(\mathcal{E}) \subseteq \mathcal{E}(\mathcal{E}_1, \mathcal{A}_i, \mathcal{E}_2)$ for some i . Since the converse inclusion is obvious, $\mathcal{F}(\mathcal{E}) = \mathcal{E}(\mathcal{E}_1, \mathcal{A}_i, \mathcal{E}_2)$, a contradiction of Lemma 1.10.

1.11 LEMMA. For $m \geq 1$, $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$ iff either

(a) $\mathcal{P}_1 \subseteq \mathcal{F}'_1, \dots, \mathcal{P}_{i-1} \subseteq \mathcal{F}'_1, \mathcal{P}_i \odot \cdots \odot \mathcal{P}_{j-1} \subseteq \mathcal{A}', \mathcal{P}_j \subseteq \mathcal{F}'_2, \dots, \mathcal{P}_m \subseteq \mathcal{F}'_2$ for some i and j , $1 \leq i \leq j \leq m+1$, or

(b) for some i , $1 \leq i \leq m$, $\mathcal{F}_1, \mathcal{A}$, and \mathcal{F}_2 , $\mathcal{P}_i = \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, and $\mathcal{P}_1 \subseteq \mathcal{F}'_1, \dots, \mathcal{P}_{i-1} \subseteq \mathcal{F}'_1, \mathcal{F}_1 \subseteq \mathcal{F}'_1, \mathcal{A} \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2), \mathcal{F}_2 \subseteq \mathcal{F}'_2, \mathcal{P}_{i+1} \subseteq \mathcal{F}'_2, \dots, \mathcal{P}_m \subseteq \mathcal{F}'_2$.

Proof. If (a) holds then, since \mathcal{F}'_1 and \mathcal{F}'_2 are closed under product,

$$\begin{aligned}(\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_{i-1}) \odot (\mathcal{P}_i \odot \cdots \odot \mathcal{P}_{j-1}) \odot (\mathcal{P}_j \odot \cdots \odot \mathcal{P}_m) &\subseteq \mathcal{F}'_1 \odot \mathcal{A}' \odot \mathcal{F}'_2 \\ &\subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2).\end{aligned}$$

Suppose (b) holds. Then

$$\mathcal{P}_i = \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2), \mathcal{F}'_2) = \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2),$$

so

$$\begin{aligned}(\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_{i-1}) \odot \mathcal{P}_i \odot (\mathcal{P}_{i+1} \odot \cdots \odot \mathcal{P}_m) &\subseteq \mathcal{F}'_1 \odot \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2) \odot \mathcal{F}'_2 \\ &= \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2).\end{aligned}$$

Conversely, suppose $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$. Two alternatives arise.

Case 1. Suppose $\mathcal{P}_1 \not\subseteq \mathcal{F}'_1$ and $\mathcal{P}_m \not\subseteq \mathcal{F}'_2$. First assume $m=1$, so $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m = \mathcal{P}_1$. Since \mathcal{P}_1 is prime, it has the form $\mathcal{P}_1 = \mathcal{F}(\mathcal{E})$ or $\mathcal{P}_1 = \mathcal{E}(\mathcal{E}_1, \mathcal{A}, \mathcal{E}_2)$ by [7, Corollary 3 of the prime decomposition theorem]. In the former case,

$$\mathcal{P}_1 \odot \mathcal{P}_1 = \mathcal{P}_1 = \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2).$$

By Lemma 1.9, $\mathcal{P}_1 = \mathcal{P}_1 \odot \mathcal{P}_1 \subseteq \mathcal{F}'_1 \odot \mathcal{O}' \odot \mathcal{F}'_2$. Since \mathcal{P}_1 is prime, it follows that $\mathcal{P}_1 \subseteq \mathcal{O}'$. Thus (a) holds for $i = 1$ and $j = 2$. In the latter case, $\mathcal{P}_1 = \mathcal{E}(\mathcal{F}_1, \mathcal{O}, \mathcal{F}_2)$, where each $\mathcal{F}_i = \mathcal{F}(\mathcal{G}_i)$ by Lemma 1.2. Then $\mathcal{P}_1 = \mathcal{F}_1 \odot \mathcal{P}_1 \odot \mathcal{F}_2$ by Corollary 1.5, so $\mathcal{P}_1 = \mathcal{F}_1 \odot \mathcal{P}_1 = \mathcal{P}_1 \odot \mathcal{F}_2$. Thus, $\mathcal{F}_1 \odot \mathcal{P}_1 = \mathcal{P}_1 \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{O}', \mathcal{F}'_2)$ and $\mathcal{P}_1 \subseteq \mathcal{F}'_2$. Since \mathcal{F}_1 and \mathcal{P}_1 are prime, they are additively prime. By Lemma 1.9, if $\mathcal{F}_1 \not\subseteq \mathcal{F}'_1$, then $\mathcal{P}_1 = \mathcal{F}_1 \odot \mathcal{P}_1 \subseteq \mathcal{F}'_1 \odot \mathcal{O}' \odot \mathcal{F}'_2$, whence $\mathcal{P}_1 \subseteq \mathcal{O}'$ and (a) holds. Similarly, if $\mathcal{F}_2 \not\subseteq \mathcal{F}'_2$, then $\mathcal{P}_1 \subseteq \mathcal{O}'$ and (a) holds. On the other hand, if $\mathcal{F}_1 \subseteq \mathcal{F}'_1$ and $\mathcal{F}_2 \subseteq \mathcal{F}'_2$, then (b) holds (since $\mathcal{O} \subseteq \mathcal{E}(\mathcal{F}_1, \mathcal{O}, \mathcal{F}_2) = \mathcal{P}_1 \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{O}', \mathcal{F}'_2)$).

Now assume $m \geq 2$. Then

$$\mathcal{P}_1 \odot (\mathcal{P}_2 \odot \cdots \odot \mathcal{P}_m) \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{O}', \mathcal{F}'_2),$$

where $\mathcal{P}_1 \not\subseteq \mathcal{F}'_1$ and $\mathcal{P}_2 \odot \cdots \odot \mathcal{P}_m = \mathcal{P}_m \not\subseteq \mathcal{F}'_2$. By Lemma 1.9, $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{F}'_1 \odot \mathcal{O}' \odot \mathcal{F}'_2$. Since each \mathcal{P}_i is prime, [7, Proposition E] implies that

$$\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_{i-1} \subseteq \mathcal{F}'_1, \quad \mathcal{P}_i \odot \cdots \odot \mathcal{P}_{j-1} \subseteq \mathcal{O}', \quad \mathcal{P}_j \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{F}'_2$$

for some i and j , $1 \leq i \leq j \leq m + 1$. Hence (a) is satisfied.

Case 2. Suppose $\mathcal{P}_1 \subseteq \mathcal{F}'_1$ or $\mathcal{P}_m \subseteq \mathcal{F}'_2$. Let p be the least index for which $\mathcal{P}_p \not\subseteq \mathcal{F}'_1$ and q the greatest index for which $\mathcal{P}_q \not\subseteq \mathcal{F}'_2$. If p does not exist, then $\mathcal{P}_i \subseteq \mathcal{F}'_1$ for all i . If q does not exist, then $\mathcal{P}_i \subseteq \mathcal{F}'_2$ for all i . If p, q exist and $p > q$, then (i) $\mathcal{P}_i \subseteq \mathcal{F}'_1$ for all i , $1 \leq i < p$, and (ii) $\mathcal{P}_i \subseteq \mathcal{F}'_2$, for all i , $q < p \leq i \leq m$. In all cases, (a) is satisfied. Suppose p and q exist, with $p \leq q$. Then

$$\begin{aligned} \mathcal{P}_1 \subseteq \mathcal{F}'_1, \dots, \mathcal{P}_{p-1} \subseteq \mathcal{F}'_1, \quad \mathcal{P}_p \odot \cdots \odot \mathcal{P}_q \subseteq \mathcal{E}(\mathcal{F}'_1, \mathcal{O}', \mathcal{F}'_2), \\ \mathcal{P}_{q+1} \subseteq \mathcal{F}'_2, \dots, \mathcal{P}_m \subseteq \mathcal{F}'_2. \end{aligned}$$

Since $\mathcal{P}_p \not\subseteq \mathcal{F}'_1$ and $\mathcal{P}_q \not\subseteq \mathcal{F}'_2$, Case 1 applies to $\mathcal{P}_p \odot \cdots \odot \mathcal{P}_q$. Then (a) or (b) holds for $\mathcal{P}_p \odot \cdots \odot \mathcal{P}_q$, from which it follows that (a) or (b) also holds for $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$. ■

2. DECISION PROCEDURE FOR EQUALITY

Using the technical results of the previous section, we now turn our attention to establishing the existence of a decision procedure for equality of grammatical families (Theorem 2.4). As indicated in the Introduction, we accomplish this by first giving a decision procedure for containment (Theorem 2.3).

Each proper grammatical family \mathcal{G} can be built up from \mathcal{P} by a finite sequence of the operations \mathcal{E} , \mathcal{F} , \odot , and \oplus [4]. We shall therefore define formal expressions which depict the way proper grammatical families are constructed by these operations. We shall extend our notation to allow \mathcal{E} to denote such an expression as well as the underlying grammatical family, permitting the context to make clear whether an expression or a family is meant. Similarly, we shall define \mathcal{O} -expressions,

\mathcal{P} -expressions, \mathcal{F} -expressions, and \mathcal{E} -expressions, and shall extend our earlier notation to permit the symbols \mathcal{A} , \mathcal{P} , \mathcal{F} , and \mathcal{E} to denote formal expressions of these four types as well as their underlying families. In the following recursive definition, universal quantification is understood. For example, " $\mathcal{E}(\mathcal{F}, \mathcal{A}, \mathcal{F}')$ is a \mathcal{E} -expression" means that, for all \mathcal{F} -expressions \mathcal{F} and \mathcal{F}' , and \mathcal{A} -expressions \mathcal{A} , the formal expressions $\mathcal{E}(\mathcal{F}, \mathcal{A}, \mathcal{F}')$ is a \mathcal{E} -expression.

DEFINITION. The collections of \mathcal{E} -expressions, \mathcal{A} -expressions, \mathcal{P} -expressions, \mathcal{F} -expressions, and \mathcal{E} -expressions are the smallest collections of formal expressions over the 8-symbol alphabet

$$\mathcal{R} \quad \hat{\mathcal{F}} \quad \mathcal{E} \quad \odot \quad \oplus \quad , \quad (\quad)$$

satisfying the following conditions:

- (1) \mathcal{R} is an \mathcal{F} -expression and a \mathcal{P} -expression,
- (2) $\hat{\mathcal{F}}(\mathcal{E})$ is an \mathcal{F} -expression and a \mathcal{P} -expression,
- (3) $\mathcal{E}(\mathcal{F}, \mathcal{A}, \mathcal{F}')$ is a \mathcal{E} -expression and a \mathcal{P} -expression,
- (4) $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$ is an \mathcal{A} -expression for $n \geq 1$,
- (5) $\mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_n$ is a \mathcal{E} -expression for $n \geq 1$.

Note that the set of \mathcal{P} -expressions is just the union of the set of \mathcal{F} -expressions with the set of \mathcal{E} -expressions. In addition, every \mathcal{P} -expression is an \mathcal{A} -expression, and every \mathcal{A} -expression is a \mathcal{E} -expression.

DEFINITION. The *value* of a grammatical expression is the family obtained by identifying \mathcal{R} with the family of regular sets, and $\hat{\mathcal{F}}$, \mathcal{E} , \odot , and \oplus with the corresponding operations on families. Recall that the operations \odot and \oplus are associative. By convention, \odot has a higher binding power than \oplus .)

The value of a \mathcal{E} -, \mathcal{A} -, or \mathcal{P} -expression is a proper grammatical family, an additively prime proper grammatical family, or a prime proper grammatical family, respectively. (See [7, Propositions A–D].)

Each grammatical expression denotes both (a) a formal expression, which is a string of symbols, and (b) its value, which is a proper grammatical family. We shall write $\mathcal{E} = \mathcal{E}'$ to mean that the expressions \mathcal{E} and \mathcal{E}' have equal values, i.e., denote the same grammatical families, and shall write $\mathcal{E} \equiv \mathcal{E}'$ to mean that \mathcal{E} and \mathcal{E}' are identical formal expressions. Thus, $\hat{\mathcal{F}}(\mathcal{R}) = \mathcal{R}$ is true, but $\hat{\mathcal{F}}(\mathcal{R}) \equiv \mathcal{R}$ is false. Similarly, we shall write $\mathcal{E} \subseteq \mathcal{E}'$ to mean that the grammatical family denoted by the first expression is a subfamily of that denoted by the second.

CONVENTION. All expressions of the form $\mathcal{A}_{i_1} \oplus \cdots \oplus \mathcal{A}_{i_n}$, where (i_1, \dots, i_n) is a permutation of $(1, \dots, n)$, are considered identical expressions.

Thus, $\hat{\mathcal{F}}(\mathcal{R}) \oplus \mathcal{R} \equiv \mathcal{R} \oplus \hat{\mathcal{F}}(\mathcal{R})$.

Note that each expression has the form $\mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_m$, $m \geq 1$, and hence the form

$$(\mathcal{P}_{11} \odot \cdots \odot \mathcal{P}_{1n_1}) \oplus \cdots \oplus (\mathcal{P}_{m1} \odot \cdots \odot \mathcal{P}_{mn_m}),$$

each $n_i \geq 1$. (The grouping parentheses, included for clarity, and the symbols "... between \odot and \odot , and \oplus and \oplus , are metasymbols and do not occur in the actual formal expression.) Thus, a grammatical expression for a proper grammatical family represents a decomposition into a sum of products of prime families (where, in addition, each prime family might be further decomposed into smaller families under \mathcal{F} and \mathcal{E}). However, this decomposition is not necessarily the unique minimal representation of the family in terms of primes.

We now show that each proper grammatical family is represented by an expression.

2.1 LEMMA. *Every proper grammatical family has an effectively calculable grammatical expression.*

Proof. If we had defined grammatical expressions by permitting arbitrary applications of the operation symbols \mathcal{F} , \mathcal{E} , \oplus , and \odot to \mathcal{R} , then Lemma 2.1 would be immediate from [4]. However, we have only allowed $\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$ to be an expression if \mathcal{G}_1 and \mathcal{G}_3 are \mathcal{F} -expressions and \mathcal{G}_2 is an \mathcal{A} -expression; we have only allowed $\mathcal{G}_1 \odot \cdots \odot \mathcal{G}_n$, $n \geq 2$, to be an expression if each \mathcal{G}_i is a \mathcal{P} -expression; and we have only allowed $\mathcal{G}_1 \oplus \cdots \oplus \mathcal{G}_n$, $n \geq 2$, to be an expression if each \mathcal{G}_i is an \mathcal{A} -expression. That these restrictions are of no consequence follows from the identity (see Lemmas 1.1 and 1.2)

$$\begin{aligned} & \mathcal{E}(\mathcal{G}', \mathcal{G}_1 \oplus \cdots \oplus \mathcal{G}_n, \mathcal{G}'') \\ &= \mathcal{E}(\mathcal{F}(\mathcal{G}'), \mathcal{G}_1, \mathcal{F}(\mathcal{G}'')) \oplus \cdots \oplus \mathcal{E}(\mathcal{F}(\mathcal{G}'), \mathcal{G}_n, \mathcal{F}(\mathcal{G}')), \end{aligned}$$

together with the facts that \oplus and \odot are associative and \odot distributes over \oplus . ■

Lemma 2.2 is the key to the decision procedure for determining whether one grammatical family includes another.

2.2 INCLUSION LEMMA. *Let \mathcal{G} and \mathcal{G}' be grammatical expressions. Then $\mathcal{G} \subseteq \mathcal{G}'$ iff one of the following holds:*

(1) $\mathcal{G} \equiv \mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_m$ and $\mathcal{G}' \equiv \mathcal{A}'_1 \oplus \cdots \oplus \mathcal{A}'_n$, with $m \geq 2$ or $n \geq 2$; and for each i there is a j such that $\mathcal{A}_i \subseteq \mathcal{A}'_j$.

(2) $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$, with $m \geq 1$; $\mathcal{G}' \equiv \mathcal{P}'_1 \odot \cdots \odot \mathcal{P}'_n$, with $n \geq 2$; and there exists a sequence $1 = i_0 \leq i_1 \leq \cdots \leq i_n = m + 1$ such that

$$\mathcal{P}_{i_{j-1}} \odot \cdots \odot \mathcal{P}_{i_j-1} \subseteq \mathcal{P}'_j \quad \text{for all } j, 1 \leq j \leq n, \text{ for which } i_{j-1} < i_j.$$

- (3) $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$, with $m \geq 1$; $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$; and either
 (a) there is an i and j , $1 \leq i \leq j \leq m+1$, such that

$$\begin{aligned} \mathcal{P}_1 \subseteq \mathcal{F}'_1, \dots, \mathcal{P}_{i-1} \subseteq \mathcal{F}'_1; \quad \mathcal{P}_i \odot \cdots \odot \mathcal{P}_{j-1} \subseteq \mathcal{A}' \quad \text{if } i < j; \\ \mathcal{P}_j \subseteq \mathcal{F}'_2, \dots, \mathcal{P}_m \subseteq \mathcal{F}'_2; \end{aligned}$$

or

- (b) there is an i , $1 \leq i \leq m$, such that $\mathcal{P}_i = \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$ and

$$\begin{aligned} \mathcal{P}_1 \subseteq \mathcal{F}'_1, \dots, \mathcal{P}_{i-1} \subseteq \mathcal{F}'_1, \quad \mathcal{F}_1 \subseteq \mathcal{F}'_1; \quad \mathcal{A} \subseteq \mathcal{G}'; \quad \mathcal{F}_2 \subseteq \mathcal{F}'_2, \\ \mathcal{P}_{i+1} \subseteq \mathcal{F}'_2, \dots, \mathcal{P}_m \subseteq \mathcal{F}'_2. \end{aligned}$$

- (4) $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$, with $m \geq 2$; \mathcal{G}' is an \mathcal{F} -expression; and $\mathcal{P}_i \subseteq \mathcal{G}'$ for all i .

- (5) $\mathcal{G} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$; $\mathcal{G}' \equiv \mathcal{F}(\mathcal{G}_0)$; and $\mathcal{G} \subseteq \mathcal{G}'_0$.

- (6) $\mathcal{G} \equiv \mathcal{F}(\mathcal{G}_0)$; \mathcal{G}' is an \mathcal{F} -expression; and $\mathcal{G}_0 \subseteq \mathcal{G}'$.

- (7) $\mathcal{G} \equiv \mathcal{R}$.

Proof. If (1), (2), or (7) holds, then obviously $\mathcal{G} \subseteq \mathcal{G}'$. If (3) holds, then $\mathcal{G} \subseteq \mathcal{G}'$ by Lemma 1.11. If (4) holds, then $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m \subseteq \mathcal{G}'$ since full AFL are closed under union and concatenation. If (5) holds, then $\mathcal{G} \subseteq \mathcal{G}'_0 \subseteq \mathcal{F}(\mathcal{G}'_0) \equiv \mathcal{G}'$. And if (6) holds, then $\mathcal{G} \equiv \mathcal{F}(\mathcal{G}_0) \subseteq \mathcal{F}(\mathcal{G}') = \mathcal{G}'$ since \mathcal{G}' is an \mathcal{F} -expression.

Now suppose that $\mathcal{G} \subseteq \mathcal{G}'$. We shall show that one of the conditions (1)–(7) applies. Three cases arise.

Case 1. Suppose either \mathcal{G} or \mathcal{G}' is not an \mathcal{A} -expression. Then $\mathcal{G} \equiv \mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_m$ and $\mathcal{G}' \equiv \mathcal{A}'_1 \oplus \cdots \oplus \mathcal{A}'_n$, where $m \geq 2$ or $n \geq 2$. For each i , $\mathcal{A}_i \subseteq \mathcal{G} \subseteq \mathcal{G}' \equiv \mathcal{A}'_1 \oplus \cdots \oplus \mathcal{A}'_n$, so $\mathcal{A}_i \subseteq \mathcal{A}'_j$ for some j (since \mathcal{A}_i is additively prime). Hence, condition (1) is satisfied.

Case 2. Suppose \mathcal{G} and \mathcal{G}' are \mathcal{A} -expressions but \mathcal{G}' is not a \mathcal{P} -expression. Then $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$, $m \geq 1$, and $\mathcal{G}' \equiv \mathcal{P}'_1 \odot \cdots \odot \mathcal{P}'_n$, $n \geq 2$. By $n-1$ applications of [7, Proposition E], condition (2) is satisfied.

Case 3. Suppose \mathcal{G} is an \mathcal{A} -expression and \mathcal{G}' is a \mathcal{P} -expression. If \mathcal{G}' is a \mathcal{E} -expression, then condition (3) follows from Lemma 1.11. If \mathcal{G}' is not a \mathcal{E} -expression, then it is an \mathcal{F} -expression.

Now if $\mathcal{G} \equiv \mathcal{R}$, then condition (7) applies. If $\mathcal{G} \equiv \mathcal{F}(\mathcal{G}_0)$, then $\mathcal{G}_0 \subseteq \mathcal{G} \subseteq \mathcal{G}'$ and condition (6) applies. Suppose $\mathcal{G} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$. Then $\mathcal{G} \subseteq \mathcal{G}'$ implies $\mathcal{G}' \neq \mathcal{R}$, so $\mathcal{G}' \equiv \mathcal{F}(\mathcal{G}'_0)$ and condition (5) follows from Lemma 1.8. Finally, suppose \mathcal{G} has none of these forms, i.e., \mathcal{G} is not a \mathcal{P} -expression. Then $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_m$ for some $m \geq 2$ and $\mathcal{P}_i \subseteq \mathcal{G} \subseteq \mathcal{G}'$ for all i , i.e., condition (4) applies. ■

We are now ready for the result on the decidability of containment.

2.3 THEOREM (Decision Procedure for Inclusion). *There is a decision procedure for determining whether one grammatical family is included in another.*

Proof. Given a grammar form, we can algorithmically determine whether the corresponding grammatical family is either \mathcal{L}_\emptyset , \mathcal{L}_ε , \mathcal{L}_{lin} , \mathcal{L}_{CF} , or is proper [3]. Hence, given two grammar forms, we can determine whether at least one of the corresponding grammatical families is not proper. Suppose at least one is not proper. Then we can immediately determine whether or not the first is included in the second, since $\mathcal{L}_\emptyset \not\subseteq \mathcal{L}_\varepsilon \not\subseteq \mathcal{L}_{\text{lin}} \not\subseteq \mathcal{L}_{\text{CF}}$ for all proper grammatical families \mathcal{G} . On the other hand, suppose both are proper. By Lemma 2.1, we can (effectively) construct grammatical expressions \mathcal{G} and \mathcal{G}' for the two families. Then Inclusion Lemma 2.2 can be used to test whether $\mathcal{G} \subseteq \mathcal{G}'$. Indeed, each case in Inclusion Lemma 2.2 reduces the problem of determining whether $\mathcal{G}_1 \subseteq \mathcal{G}_2$ to a finite set of similar problems, each of shorter length. (The length of " $\mathcal{G}_1 \subseteq \mathcal{G}_2$ " is defined to be the sum of the lengths of the formal expressions \mathcal{G}_1 and \mathcal{G}_2 .) Hence, Inclusion Lemma 2.2 gives rise to a terminating recursive procedure for determining whether $\mathcal{G}_1 \subseteq \mathcal{G}_2$. ■

EXAMPLE. Consider the grammar forms $G_1 = (\{\sigma, \xi, a\}, \{a\}, P_1, \sigma)$ and $G_2 = (\{\sigma, \xi, a\}, \{a\}, P_2, \sigma)$, where

$$P_1 = \{\sigma \rightarrow \xi\sigma, \sigma \rightarrow \xi, \xi \rightarrow a\xi a, \xi \rightarrow \varepsilon\}$$

and

$$P_2 = \{\sigma \rightarrow a\sigma a, \sigma \rightarrow \xi\xi, \xi \rightarrow a\xi a, \xi \rightarrow \varepsilon\}.$$

Applying the construction in [4] yields grammatical expressions such as

$$\mathcal{G}_1 \equiv \hat{\mathcal{F}}(\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})) \oplus \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})$$

and

$$\mathcal{G}_2 \equiv \mathcal{E}(\mathcal{R}, \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})) \odot \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}), \mathcal{R})$$

for the grammatical families generated by G_1 and G_2 , respectively. Such expressions will in general not be minimal expressions. (The next section shows how to obtain minimal grammatical expressions. A minimal expression for \mathcal{G}_1 is $\hat{\mathcal{F}}(\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}))$, while the expression obtained for \mathcal{G}_2 is already minimal.)

To illustrate the recursive use of Inclusion Lemma 2.2, let us blindly test whether $\mathcal{G}_1 \subseteq \mathcal{G}_2$. In practice, of course, the hand computation would proceed more quickly since various simplifications would be made during the computation.

By Lemma 2.2(1), $\mathcal{G}_1 \subseteq \mathcal{G}_2$ iff $\hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \subseteq \mathcal{G}_2$ and $\mathcal{G}_{\text{lin}} \subseteq \mathcal{G}_2$, where $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})$ is abbreviated as \mathcal{G}_{lin} (since it is the family of linear languages). Then by Lemma 2.2(3a) ((3b) is not applicable),

$$\hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \subseteq \mathcal{G}_2 \equiv \mathcal{E}(\mathcal{R}, \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}}, \mathcal{R})$$

iff

$$\hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \subseteq \mathcal{R} \quad \text{or} \quad \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \subseteq \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}}.$$

But $\hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \subseteq \mathcal{R}$ iff $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}) \equiv \mathcal{G}_{\text{lin}} \subseteq \mathcal{R}$, by Lemma 2.2(6). This last inclusion is false since none of the seven conditions in Lemma 2.2 is satisfied. Now

$$\begin{aligned}
 \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) &\subseteq \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}} \\
 \text{iff } \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) &\subseteq \mathcal{G}_{\text{lin}} \equiv \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}) && \text{by Lemma 2.2(2),} \\
 \text{iff } \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) &\subseteq \mathcal{R} && \text{by Lemma 2.2(3a),} \\
 \text{iff } \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}) &\equiv \mathcal{G}_{\text{lin}} \subseteq \mathcal{R}, && \text{by Lemma 2.2(6).}
 \end{aligned}$$

The last inclusion is false since it does not satisfy any of the conditions in Lemma 2.2. Hence, $\mathcal{G}_1 \subseteq \mathcal{G}_2$ is false.

Similarly, we can test whether $\mathcal{G}_2 \subseteq \mathcal{G}_1$. This time, however, let us take some shortcuts. Obviously $\mathcal{G}_1 \equiv \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \oplus \mathcal{G}_{\text{lin}} = \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}})$, so

$$\begin{aligned}
 \mathcal{E}(\mathcal{R}, \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}}, \mathcal{R}) &\equiv \mathcal{G}_2 \subseteq \mathcal{G}_1 = \hat{\mathcal{F}}(\mathcal{G}_{\text{lin}}) \\
 \text{iff } \mathcal{E}(\mathcal{R}, \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}}, \mathcal{R}) &\subseteq \mathcal{G}_{\text{lin}} = \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}) && \text{by Lemma 2.2(5),} \\
 \text{iff } \mathcal{E}(\mathcal{R}, \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}}, \mathcal{R}) &\subseteq \mathcal{R} && \text{by Lemma 2.2(3a),} \\
 \text{or } \mathcal{G}_{\text{lin}} \odot \mathcal{G}_{\text{lin}} &\subseteq \mathcal{G}_{\text{lin}} && \text{by Lemma 2.2(3b).}
 \end{aligned}$$

But the latter containments are clearly false. Thus, \mathcal{G}_1 and \mathcal{G}_2 are incomparable. ■

Since $\mathcal{G}_1 = \mathcal{G}_2$ iff $\mathcal{G}_1 \subseteq \mathcal{G}_2$, and $\mathcal{G}_2 \subseteq \mathcal{G}_1$, Theorem 2.3 allows us to state the existence of a decision procedure for equality for grammatical families. Indeed, it was the question of whether such a procedure existed that initiated [7] and the present investigation.

2.4 THEOREM. *There is a decision procedure for determining whether two grammatical families are equal.* ■

The following two results are also corollaries of Theorem 2.3:

2.5 COROLLARY. *The unique representation of a grammatical family as a minimal sum of products of prime grammatical families can be effectively calculated.*

Proof. By Lemma 2.1, a representation of the given family as a sum of products of primes can be effectively constructed. This representation can then be reduced to a minimal representation using Theorem 2.3. ■

2.6 COROLLARY. *There is a decision procedure for determining whether a grammatical family is (a) prime or (b) additively prime.*

Proof. This follows from Corollary 2.5 and the fact ([7, Corollary 1 of the prime decomposition theorem]) that a grammatical family is prime (additively prime) if and only if its unique representation as a minimal sum of products of primes consists of just a single factor (summand). ■

3. CANONICAL EXPRESSIONS

By the prime decomposition theorem, every grammatical family can be uniquely decomposed into a minimal sum of products of prime grammatical families. If the prime families could in turn be uniquely decomposed into smaller families, then this process could be continued and the prime decomposition theorem could be extended so that the given grammatical family is uniquely analyzed into a combination of such basic "atomic" families as the family \mathcal{R} of regular languages. Unfortunately, the representation of proper grammatical families in terms of \mathcal{R} in the preceding section (Lemma 2.1) is not unique. For example, the following pairs of expressions represent the same families but are not identical (for all admissible values of \mathcal{G} , \mathcal{A} , \mathcal{P}_1 , and \mathcal{F}_i):

$$\begin{aligned}\mathcal{A} \oplus \mathcal{A} &= \mathcal{A}, \\ \mathcal{F} \odot \mathcal{F} &= \mathcal{F}, \\ \mathcal{F}_1 \odot \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) &= \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}), \\ \mathcal{E}(\mathcal{F}_1, \mathcal{F}_1 \odot \mathcal{P}, \mathcal{F}_2) &= \mathcal{E}(\mathcal{F}_1, \mathcal{P}, \mathcal{F}_2), \\ \mathcal{E}(\mathcal{F}_1, \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2), \mathcal{F}_2) &= \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2), \\ \mathcal{F}(\mathcal{F}) &= \mathcal{F},\end{aligned}$$

and

$$\mathcal{F}(\mathcal{P}_1 \odot \mathcal{P}_2) = \mathcal{F}(\mathcal{P}_1 \oplus \mathcal{P}_2).$$

(Note that $\mathcal{A}_1 \oplus \mathcal{A}_2$ and $\mathcal{A}_2 \oplus \mathcal{A}_1$ are identical by convention, and thus not an example of equal but nonidentical expressions.) A restricted class of grammatical expressions, called the *canonical* expressions, in which redundancies such as the preceding ones cannot occur, is defined. In Theorem 3.1, it is shown that each proper grammatical family has a unique canonical expression. (Thus, two proper grammatical families are equal iff they have the same canonical expression.) Since a family's canonical expression is a refinement of its unique minimal prime decomposition, this result is the desired extension of the prime decomposition theorem.

Formally, our notion of canonical expression is given by

DEFINITION. A grammatical expression \mathcal{G} is *canonical* if it has one of the following forms:

- (1) $\mathcal{G} \equiv \mathcal{R}$.
- (2) $\mathcal{G} \equiv \mathcal{F}(\mathcal{E}_1 \oplus \dots \oplus \mathcal{E}_n)$, where $n \geq 1$ and $\mathcal{E}_1 \oplus \dots \oplus \mathcal{E}_n$ is canonical.
- (3) $\mathcal{G} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, where $\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2$ are canonical and
 - (a) if $\mathcal{A} \equiv \mathcal{P}_1 \odot \dots \odot \mathcal{P}_n \neq \mathcal{R}$, with $n \geq 1$, then $\mathcal{P}_1 \notin \mathcal{F}_1$ and $\mathcal{P}_n \notin \mathcal{F}_2$,

and

(b) if $\mathcal{A} \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$, then $(\mathcal{F}'_1, \mathcal{F}'_2)$ and $(\mathcal{F}_1, \mathcal{F}_2)$ are incomparable.

(4) $\mathcal{G} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$, where $n \geq 2$, each \mathcal{P}_i is canonical, and the product is minimal.

(5) $\mathcal{G} \equiv \mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_n$, where $n \geq 2$, each \mathcal{A}_i is canonical, and the sum is minimal.

Our result on canonical expressions is

3.1 THEOREM (Canonical decomposition theorem). *Every proper grammatical family \mathcal{L} has a unique canonical expression*

$$(\mathcal{P}_{11} \odot \cdots \odot \mathcal{P}_{1n_1}) \oplus \cdots \oplus (\mathcal{P}_{m1} \odot \cdots \odot \mathcal{P}_{mn_m}),$$

where $m \geq 1$ and $n_i \geq 1$ for each i ; and this expression is effectively calculable. Furthermore, if \mathcal{L}_{ij} is the value of \mathcal{P}_{ij} for each i and j , then

$$(\mathcal{L}_{11} \odot \cdots \odot \mathcal{L}_{1n_1}) \oplus \cdots \oplus (\mathcal{L}_{m1} \odot \cdots \odot \mathcal{L}_{mn_m})$$

is the unique minimal decomposition of \mathcal{L} into a sum of products of prime families.

Proof. By Lemma 2.1, \mathcal{L} has an effectively calculable grammatical expression \mathcal{G} . The following effective recursive procedure converts \mathcal{G} into a canonical expression without changing its type (i.e., \mathcal{F} -, \mathcal{E} -, and \mathcal{A} -expressions remain \mathcal{F} -, \mathcal{E} -, and \mathcal{A} -expressions):

- (1) If $\mathcal{G} \equiv \mathcal{R}$, then \mathcal{G} is already canonical.
- (2) If $\mathcal{G} \equiv \hat{\mathcal{F}}(\mathcal{G}')$, then convert \mathcal{G}' into a canonical expression

$$(\mathcal{P}_{11} \odot \cdots \odot \mathcal{P}_{1n_1}) \oplus \cdots \oplus (\mathcal{P}_{m1} \odot \cdots \odot \mathcal{P}_{mn_m}),$$

where $m \geq 1$, $n_i \geq 1$ for each i , and each \mathcal{P}_{ij} is a canonical \mathcal{P} -expression. Now change this expression to

$$\mathcal{P}_{11} \oplus \cdots \oplus \mathcal{P}_{1n_1} \oplus \cdots \oplus \mathcal{P}_{m1} \oplus \cdots \oplus \mathcal{P}_{mn_m}.$$

Then delete each \mathcal{P}_{ij} of the form \mathcal{R} and replace each \mathcal{P}_{ij} of the form $\hat{\mathcal{F}}(\mathcal{G}_{ij})$ by \mathcal{G}_{ij} . Since each \mathcal{P}_{ij} is canonical, each \mathcal{G}_{ij} is a sum of canonical \mathcal{E} -expressions. Furthermore, each \mathcal{P}_{ij} not of the form \mathcal{R} or $\hat{\mathcal{F}}(\mathcal{G}_{ij})$ and hence surviving unchanged is already a canonical \mathcal{E} -expression. Thus, \mathcal{G}' has been converted into a sum of zero or more canonical \mathcal{E} -expressions. Now use the decision procedure for inclusion to remove \mathcal{E} -expressions from this sum until it is minimal. This converts \mathcal{G}' into \mathcal{L}' , where \mathcal{L}' is empty or is a minimal sum of canonical \mathcal{E} -expressions. In the former case, replace \mathcal{G} by \mathcal{R} ; in the latter case, by $\hat{\mathcal{F}}(\mathcal{L}')$. That this process does not

change the value of \mathcal{G} follows from associativity of \oplus and \odot , commutativity of \oplus , and repeated applications of the easily verified identities:

$$\hat{\mathcal{F}}(\mathcal{G}_1 \oplus (\mathcal{G}_2 \odot \mathcal{G}_3)) = \hat{\mathcal{F}}(\mathcal{G}_1 \oplus \mathcal{G}_2 \oplus \mathcal{G}_3),$$

$$\hat{\mathcal{F}}(\mathcal{G} \oplus \mathcal{R}) = \hat{\mathcal{F}}(\mathcal{G}),$$

$$\hat{\mathcal{F}}(\mathcal{G} \oplus \hat{\mathcal{F}}(\mathcal{G}_2)) = \hat{\mathcal{F}}(\mathcal{G}_1 \oplus \mathcal{G}_2),$$

$$\mathcal{G}_1 \oplus \mathcal{G}_2 = \mathcal{G}_1 \quad \text{if } \mathcal{G}_2 \subseteq \mathcal{G}_1,$$

and

$$\hat{\mathcal{F}}(\phi) = \mathcal{R}.$$

(3) If $\mathcal{G} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, then convert \mathcal{F}_1 , \mathcal{A} , and \mathcal{F}_2 to canonical expressions \mathcal{F}'_1 , \mathcal{A}' , and \mathcal{F}'_2 . Let $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$. Keep performing the following steps as long as any is applicable.

(a) If $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{P}_1 \odot \dots \odot \mathcal{P}_n, \mathcal{F}'_2)$, where $n \geq 2$ and $\mathcal{P}_1 \subseteq \mathcal{F}'_1$, then delete \mathcal{P}_1 from \mathcal{G}' .

(b) If $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{P}_1 \odot \dots \odot \mathcal{P}_n, \mathcal{F}'_2)$, where $n \geq 2$ and $\mathcal{P}_n \subseteq \mathcal{F}'_2$, then delete \mathcal{P}_n from \mathcal{G}' .

(c) If $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{P}, \mathcal{F}'_2)$, where $\mathcal{P} \neq \mathcal{R}$ and either $\mathcal{P} \subseteq \mathcal{F}'_1$ or $\mathcal{P} \subseteq \mathcal{F}'_2$, then replace \mathcal{P} by \mathcal{R} in \mathcal{G}' .

(d) If $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{E}(\mathcal{F}''_1, \mathcal{A}'', \mathcal{F}''_2), \mathcal{F}'_2)$ and $(\mathcal{F}''_1, \mathcal{F}''_2) \subseteq (\mathcal{F}'_1, \mathcal{F}'_2)$, then replace \mathcal{G}' by $\mathcal{E}(\mathcal{F}'_1, \mathcal{A}'', \mathcal{F}'_2)$.

(e) If $\mathcal{G}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{E}(\mathcal{F}''_1, \mathcal{A}'', \mathcal{F}''_2), \mathcal{F}'_2)$ and $(\mathcal{F}'_1, \mathcal{F}'_2) \subseteq (\mathcal{F}''_1, \mathcal{F}''_2)$, then replace \mathcal{G}' by $\mathcal{E}(\mathcal{F}''_1, \mathcal{A}'', \mathcal{F}''_2)$.

Since each application of (a), (b), (d), or (e) shortens \mathcal{G}' while an application of (c) cannot be followed by (a)–(e), this procedure eventually terminates. The end result is a \mathcal{E} -expression which must be canonical since (a)–(e) are no longer applicable. The process does not change the value of \mathcal{G}' because of the following identities:

$$\mathcal{E}(\mathcal{G}_1, \mathcal{G}'_1 \odot \mathcal{G}_2, \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \quad \text{if } \mathcal{G}'_1 \subseteq \mathcal{G}_1 \text{ by Corollary 1.7,}$$

$$\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2 \odot \mathcal{G}'_3, \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) \quad \text{if } \mathcal{G}'_3 \subseteq \mathcal{G}_3 \text{ by Corollary 1.7,}$$

$$\mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2 \odot \mathcal{R}, \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{R} \odot \mathcal{G}_2, \mathcal{G}_3)$$

$$= \mathcal{E}(\mathcal{G}_1, \mathcal{R}, \mathcal{G}_3)$$

$$\text{if } \mathcal{G}_2 \subseteq \mathcal{G}_1 \text{ or } \mathcal{G}_2 \subseteq \mathcal{G}_3 \text{ by the preceding identities,}$$

$$\mathcal{E}(\mathcal{G}_1, \mathcal{E}(\mathcal{G}'_1, \mathcal{G}_2, \mathcal{G}'_3), \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$$

$$\text{if } (\mathcal{G}'_1, \mathcal{G}'_3) \subseteq (\mathcal{G}_1, \mathcal{G}_3) \text{ by Corollary 1.6,}$$

and

$$\mathcal{E}(\mathcal{G}_1, \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}'_3), \mathcal{G}_3) = \mathcal{E}(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3)$$

$$\text{if } (\mathcal{G}_1, \mathcal{G}_3) \subseteq (\mathcal{G}'_1, \mathcal{G}'_3) \text{ by Corollary 1.6.}$$

(4) If $\mathcal{E} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$, where $n \geq 2$, then convert each \mathcal{P}_i to a canonical \mathcal{P} -expression and delete terms until the product is minimal, hence canonical.

(5) If $\mathcal{E} \equiv \mathcal{A}_1 \oplus \cdots \oplus \mathcal{A}_n$, where $n \geq 2$, then convert each \mathcal{A}_i to a canonical \mathcal{A} -expression and delete summands until the sum is minimal, hence canonical.

It follows by induction on the length of \mathcal{E} that the procedure embodied in (1)–(5) always terminates. Since there is a decision procedure for determining the inclusion of one grammatical family in another, the procedure in (1)–(5) is effective.

Now every grammatical expression has the form

$$(\mathcal{P}_{11} \odot \cdots \odot \mathcal{P}_{1n_1}) \oplus \cdots \oplus (\mathcal{P}_{m1} \odot \cdots \odot \mathcal{P}_{mn_m}),$$

where $m \geq 1$ and $n_i \geq 1$ for each i . If the expression is canonical, then the sum and all of the products are minimal. It follows from [7, Corollaries 5 and 6 of the prime decomposition theorem] that this expression refines the minimal prime decomposition of the corresponding family.

Finally, it remains to show that if \mathcal{E} and \mathcal{E}' are canonical expressions for the same family, then $\mathcal{E} \equiv \mathcal{E}'$. We proceed by induction on the sum of the lengths of the two expressions. (The case when the sum is one is satisfied vacuously.) Since both expressions refine the unique minimal prime decomposition of the underlying family, they have the form

$$\mathcal{E} \equiv (\mathcal{P}_{11} \odot \cdots \odot \mathcal{P}_{1n_1}) \oplus \cdots \oplus (\mathcal{P}_{m1} \odot \cdots \odot \mathcal{P}_{mn_m})$$

and

$$\mathcal{E}' \equiv (\mathcal{P}'_{11} \odot \cdots \odot \mathcal{P}'_{1n_1}) \oplus \cdots \oplus (\mathcal{P}'_{m1} \odot \cdots \odot \mathcal{P}'_{mn_m}),$$

where $\mathcal{P}_{ij} = \mathcal{P}'_{ij}$ for all i, j . It therefore suffices to prove $\mathcal{P}_{ij} \equiv \mathcal{P}'_{ij}$. To establish the uniqueness of canonical expressions it thus suffices to show that equal canonical \mathcal{P} -expressions are identical. By Lemma 1.10, an \mathcal{F} -expression cannot equal a \mathcal{E} -expression. Hence, it suffices to prove that

- (1) equal canonical \mathcal{F} -expressions are identical, and
- (2) equal canonical \mathcal{E} -expressions are identical.

Consider (1). Suppose that \mathcal{F} and \mathcal{F}' are equal canonical \mathcal{F} -expressions. If $\mathcal{F} = \mathcal{F}' = \mathcal{R}$, then clearly $\mathcal{F} \equiv \mathcal{R}$ and $\mathcal{F}' \equiv \mathcal{R}$. Therefore, we may assume that $\mathcal{F}' = \mathcal{F} \neq \mathcal{R}$. Thus, $\mathcal{F} \equiv \mathcal{F}(\mathcal{E}_1 \oplus \cdots \oplus \mathcal{E}_m)$ and $\mathcal{F}' \equiv \mathcal{F}(\mathcal{E}'_1 \oplus \cdots \oplus \mathcal{E}'_n)$ for some canonical expressions $\mathcal{E}_1 \oplus \cdots \oplus \mathcal{E}_m$ and $\mathcal{E}'_1 \oplus \cdots \oplus \mathcal{E}'_n$. Since $\mathcal{F} \subseteq \mathcal{F}'$, it follows essentially from Lemma 2.2(6) and (5) that $\mathcal{E}_i \subseteq \mathcal{E}'_1 \oplus \cdots \oplus \mathcal{E}'_n$, $1 \leq i \leq m$, so $\mathcal{E}_1 \oplus \cdots \oplus \mathcal{E}_m \subseteq \mathcal{E}'_1 \oplus \cdots \oplus \mathcal{E}'_n$. The reverse inclusion is similar, so $\mathcal{E}_1 \oplus \cdots \oplus \mathcal{E}_m = \mathcal{E}'_1 \oplus \cdots \oplus \mathcal{E}'_n$. Since these canonical expressions are shorter than \mathcal{F} and \mathcal{F}' , they are identical by induction. Hence, \mathcal{F} and \mathcal{F}' are identical.

Consider (2). Suppose that $\mathcal{E} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$ and $\mathcal{E}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{A}', \mathcal{F}'_2)$ are equal canonical \mathcal{E} -expressions. Let $i = 1, 2$. By Lemma 1.10, $\mathcal{E} \neq \mathcal{F}'_i$. Since $\mathcal{F}'_i \subseteq$

$\mathcal{E}' \subseteq \mathcal{E}$, $\mathcal{E} \not\subseteq \mathcal{F}'_1$. Since $\mathcal{E} \subseteq \mathcal{E}'$, it therefore follows from Lemma 2.2(3) that either $\mathcal{E} \subseteq \mathcal{O}'$ or $(\mathcal{F}_1, \mathcal{F}_2) \subseteq (\mathcal{F}'_1, \mathcal{F}'_2)$. In the former case, $\mathcal{E} = \mathcal{O}'$ since $\mathcal{O}' \subseteq \mathcal{E}' = \mathcal{E}$. By (3) of the definition of canonical expression, \mathcal{O}' is canonical. Since the length of \mathcal{O}' is less than that of \mathcal{E}' , $\mathcal{E} \equiv \mathcal{O}'$ by induction. Thus, either $\mathcal{E} \equiv \mathcal{O}'$ or $(\mathcal{F}_1, \mathcal{F}_2) \subseteq (\mathcal{F}'_1, \mathcal{F}'_2)$. Similarly, either $\mathcal{E}' \equiv \mathcal{O}$ or $(\mathcal{F}'_1, \mathcal{F}'_2) \subseteq (\mathcal{F}_1, \mathcal{F}_2)$. Suppose that $\mathcal{E} \equiv \mathcal{O}'$. Then

$$(*) \quad \mathcal{E}' \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{E}, \mathcal{F}'_2) \equiv \mathcal{E}(\mathcal{F}'_1, \mathcal{E}(\mathcal{F}_1, \mathcal{O}, \mathcal{F}_2), \mathcal{F}'_2).$$

Since \mathcal{E}' is canonical, $(\mathcal{F}'_1, \mathcal{F}'_2) \not\subseteq (\mathcal{F}_1, \mathcal{F}_2)$ by (3b) of the definition of a canonical expression. Therefore $\mathcal{E}' \equiv \mathcal{O}$, contradicting (*). Thus $\mathcal{E} \not\equiv \mathcal{O}'$. Similarly, $\mathcal{E}' \not\equiv \mathcal{O}$. Then $(\mathcal{F}_1, \mathcal{F}_2) \subseteq (\mathcal{F}'_1, \mathcal{F}'_2) \subseteq (\mathcal{F}_1, \mathcal{F}_2)$, whence $(\mathcal{F}_1, \mathcal{F}_2) = (\mathcal{F}'_1, \mathcal{F}'_2)$. By induction, $\mathcal{F}_1 \equiv \mathcal{F}'_1$ and $\mathcal{F}_2 \equiv \mathcal{F}'_2$. To prove $\mathcal{E} \equiv \mathcal{E}'$, it therefore suffices to establish that $\mathcal{O} \equiv \mathcal{O}'$. By induction, it is enough to show $\mathcal{O} \subseteq \mathcal{O}'$. We shall prove $\mathcal{O} \subseteq \mathcal{O}'$, the reverse inclusion following by symmetry.

Consider Lemma 2.2 applied to the inclusion $\mathcal{O} \subseteq \mathcal{E} = \mathcal{E}' \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{O}', \mathcal{F}_2)$. If $\mathcal{O} \equiv \mathcal{R}$, then $\mathcal{O} \subseteq \mathcal{O}'$. Suppose $\mathcal{O} \equiv \mathcal{F}(\mathcal{G})$. Since $\mathcal{E} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{O}, \mathcal{F}_2)$, $\mathcal{O} \not\subseteq \mathcal{F}_1$, and $\mathcal{O} \not\subseteq \mathcal{F}_2$ by (3a) of the definition of a canonical expression. By Lemma 2.2(3a), $\mathcal{O} \subseteq \mathcal{O}'$. Suppose $\mathcal{O} \equiv \mathcal{E}(\mathcal{F}''_1, \mathcal{O}'', \mathcal{F}''_2)$. Since $\mathcal{E} = \mathcal{E}(\mathcal{F}_1, \mathcal{E}(\mathcal{F}''_1, \mathcal{O}'', \mathcal{F}''_2), \mathcal{F}_2)$ is canonical, $\mathcal{O} \not\subseteq \mathcal{F}_1$, $\mathcal{O} \not\subseteq \mathcal{F}_2$, and $(\mathcal{F}''_1, \mathcal{F}''_2) \not\subseteq (\mathcal{F}_1, \mathcal{F}_2)$. By Lemma 2.2(3), $\mathcal{O} \subseteq \mathcal{O}'$. Finally, suppose $\mathcal{O} \equiv \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$, where $n \geq 2$. Since $\mathcal{E} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n, \mathcal{F}_2)$ is canonical, $\mathcal{P}_1 \not\subseteq \mathcal{F}_1$ and $\mathcal{P}_n \not\subseteq \mathcal{F}_2$. By Lemma 2.2(3a), $\mathcal{O} \subseteq \mathcal{O}'$. ■

Since the canonical expression refines the minimal prime decomposition of a family, we have

3.2 COROLLARY. *A proper grammatical family is prime (respectively, additively prime) iff its canonical expression is a \mathcal{P} -expression (respectively, \mathcal{A} -expression). ■*

The canonical decomposition theorem suggests that constructing canonical expressions can be useful. However, two of the steps used in the recursive construction of canonical expressions appear difficult to implement,

(*) $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$, with $n \geq 2$, is canonical if each \mathcal{P}_i is canonical and the product is minimal.

(**) $\mathcal{O}_1 \oplus \cdots \oplus \mathcal{O}_n$, with $n \geq 2$, is canonical if each \mathcal{O}_i is canonical and the sum is minimal.

In both cases, testing for minimality appears to require using the decision procedure for inclusion to test whether the original expression denotes a subfamily of any of the n families obtained by deleting a term from the original expression. However, in the case of (**), it is clear that this reduces to testing that the \mathcal{O}_i are pairwise incomparable. A similar reduction can be applied to (*), as indicated in

3.3 PROPOSITION. (a) $\mathcal{O}_1 \oplus \cdots \oplus \mathcal{O}_n$, with $n \geq 2$, is canonical iff the \mathcal{O}_i are canonical and pairwise incomparable.

(b) $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$, with $n \geq 2$, is canonical iff the \mathcal{P}_i are canonical and the following conditions hold:

(i) For all i , $1 \leq i < n$, (a) if \mathcal{P}_{i+1} is an \mathcal{F} -expression, then

$$\mathcal{P}_i \not\subseteq \mathcal{P}_{i+1},$$

and (b) if $\mathcal{P}_{i+1} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, then

$$\mathcal{P}_i \not\subseteq \mathcal{F}_1.$$

(ii) For all i , $1 < i \leq n$, (a) if \mathcal{P}_{i-1} is an \mathcal{F} -expression, then

$$\mathcal{P}_i \not\subseteq \mathcal{P}_{i-1},$$

and (b) if $\mathcal{P}_{i-1} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$, then

$$\mathcal{P}_i \not\subseteq \mathcal{F}_2.$$

Proof. (a) Obvious.

(b) Suppose that $\mathcal{P}_i \subseteq \mathcal{P}_{i+1}$ for some i , where \mathcal{P}_{i+1} is an \mathcal{F} -expression. Then

$$\mathcal{P}_i \odot \mathcal{P}_{i+1} \subseteq \mathcal{P}_{i+1} \odot \mathcal{P}_{i+1} \subseteq \mathcal{P}_{i+1} \subseteq \mathcal{P}_i \odot \mathcal{P}_{i+1},$$

so $\mathcal{P}_i \odot \mathcal{P}_{i+1} = \mathcal{P}_{i+1}$. Suppose that $\mathcal{P}_i \subseteq \mathcal{F}_1$ for some i , where $\mathcal{P}_{i+1} \equiv \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$. Then

$$\mathcal{P}_i \odot \mathcal{P}_{i+1} \subseteq \mathcal{F}_1 \odot \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) = \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) = \mathcal{P}_{i+1} \subseteq \mathcal{P}_i \odot \mathcal{P}_{i+1},$$

so $\mathcal{P}_i \odot \mathcal{P}_{i+1} = \mathcal{P}_{i+1}$. In either case, $\mathcal{P}_i \odot \cdots \odot \mathcal{P}_n$ is not a minimal product. So if (i) fails, then $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$ is not canonical. The case where (ii) fails is similar.

Now suppose that $\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n$ is not canonical. Then

$$\mathcal{P}_1 \odot \cdots \odot \mathcal{P}_n = \mathcal{P}_1 \odot \cdots \odot \mathcal{P}_{k-1} \odot \mathcal{P}_{k+1} \odot \cdots \odot \mathcal{P}_n$$

for some k . By Lemma 2.2, there is a sequence $1 = i_0 \leq i_1 \leq \cdots \leq i_{n-1} = n + 1$ such that

$$\mathcal{P}_{i_{j-1}} \odot \cdots \odot \mathcal{P}_{i_j} \subseteq \mathcal{P}_j$$

for all j , $1 \leq j < k$, satisfying $i_{j-1} < i_j$, and

$$\mathcal{P}_{i_{j-1}} \odot \cdots \odot \mathcal{P}_{i_j} \subseteq \mathcal{P}_{j+1}$$

for all j , $k \leq j < n$, satisfying $i_{j-1} < i_j$. Let f be the nondecreasing function from $\{1, \dots, n\}$ into $\{1, \dots, n-1\}$ defined by $f(i) = j$, where j is the unique integer such that $i_{j-1} \leq i < i_j$. Two alternatives arise.

Case 1. $f(k) \leq k-1$ (so that $\mathcal{P}_1 \odot \dots \odot \mathcal{P}_k \subseteq \mathcal{P}_1 \odot \dots \odot \mathcal{P}_{k-1}$). Let $l = \min \{j \mid f(j) \leq j-1\}$. Then $1 < l \leq k$, and $f(l-1) \leq f(l) \leq l-1 \leq f(l-1)$, the last inequality following from the minimality of l . Thus $f(l-1) = f(l) = l-1$, so $\mathcal{P}_{l-1} \odot \mathcal{P}_l \subseteq \mathcal{P}_{l-1}$. If \mathcal{P}_{l-1} is an \mathcal{F} -expression, then (ii) is violated (since $\mathcal{P}_l \subseteq \mathcal{P}_{l-1} \odot \mathcal{P}_l \subseteq \mathcal{P}_{l-1}$). Suppose $\mathcal{P}_{l-1} \in \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$. Then

$$\mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) \odot \mathcal{P}_l \subseteq \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2).$$

By Lemma 2.2(3), either

- (1) $\mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) \subseteq \mathcal{F}_1$, whence $\mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) = \mathcal{F}_1$, or
- (2) $\mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) \subseteq \mathcal{A}$, whence $\mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2) = \mathcal{A}$, or
- (3) $\mathcal{P}_l \subseteq \mathcal{F}_2$.

However, (1) and (2) contradict the uniqueness of canonical expressions. Hence $\mathcal{P}_l \subseteq \mathcal{F}_2$, and (ii) is violated.

Case 2. $f(k) \geq k$ (so that $\mathcal{P}_k \odot \dots \odot \mathcal{P}_n \subseteq \mathcal{P}_{k+1} \odot \dots \odot \mathcal{P}_n$). By an argument similar to Case 1, $\mathcal{P}_l \odot \mathcal{P}_{l+1} \subseteq \mathcal{P}_{l+1}$ for some l . This implies that (i) is violated.

In either case, (i) or (ii) is false if $\mathcal{P}_1 \odot \dots \odot \mathcal{P}_n$ is not canonical. ■

We now use the recursive definition of canonical expressions to construct a table of the first few such expressions. Clause (1) of the definition introduces the expression \mathcal{R} , which is an \mathcal{F} -expression (hence also a \mathcal{P} -, \mathcal{A} -, and \mathcal{L} -expression). Since no \mathcal{E} -expressions have been constructed as yet, clause (2) is not applicable. Clause (3) produces $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R})$. This expression, whose value is the family of linear languages, will be denoted by \mathcal{L} . Clause (4) produces $\mathcal{L} \odot \mathcal{L} = \mathcal{L}^2$, $\mathcal{L} \odot \mathcal{L} \odot \mathcal{L} = \mathcal{L}^3$, (By Proposition 3.3, these products are canonical.) Letting $\mathcal{L}^0 = \mathcal{R}$, the families so far produced are \mathcal{L}^k for each $k \geq 0$. Since these families are pairwise comparable, clause (5) does not apply, by Proposition 3.3. During the second "generation," clause (2) produces $\mathcal{F}(\mathcal{L})$, denoted \mathcal{L}^ω . To simplify the description, write $k < \omega$ for all integers k . Then clause (3) produces

$$\mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R}),$$

$2 \leq k \leq \omega$ ($k=0$ is not new and $k=1$ is ruled out by (3b)),

$$\mathcal{E}(\mathcal{L}^\omega, \mathcal{L}^k, \mathcal{R}), \quad k=0, \quad \text{i.e.,} \quad \mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R}),$$

($k=1$ is ruled out by (3b) and $2 \leq k \leq \omega$ is ruled out by (3a)),

$$\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega),$$

and

$$\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega).$$

Clause (4) produces all products of the form

$$(*) \quad \mathcal{L}^{k_0} \odot \mathcal{E}_1 \odot \mathcal{L}^{k_1} \odot \dots \odot \mathcal{E}_n \odot \mathcal{L}^{k_n},$$

subject to the following: (i) $n \geq 2$, or $n = 1$, and $k_0 + k_1 > 0$, (ii) $0 \leq k_i \leq \omega$, (iii) $k_i = 0$ means that \mathcal{L}^{k_i} is absent, (iv) \mathcal{E}_i is in

$$\{\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R}), \mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R}), \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega), \mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega) \mid 2 \leq k \leq \omega\},$$

and (v) (by Theorem 3.3) if $\mathcal{E}_i = \mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R})$ or $\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega)$, then $k_{i-1} = 0$, and if $\mathcal{E}_i = \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega)$ or $\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega)$, then $k_i = 0$. Note that

$$\begin{aligned} \mathcal{E}(\mathcal{R}, \mathcal{L}^2, \mathcal{R}) &\subseteq \mathcal{E}(\mathcal{R}, \mathcal{L}^3, \mathcal{R}) \subseteq \dots \subseteq \mathcal{E}(\mathcal{R}, \mathcal{L}^\omega, \mathcal{R}) \\ &\subseteq \mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R}) \subseteq \mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega), \end{aligned}$$

and similarly with $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega)$ in place of $\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R})$; and that these are the only inclusions among families of form (*). Also, $\mathcal{L}^i \subseteq \mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R})$ iff $i \leq k$. Hence, clause (5) produces the canonical expressions

$$\mathcal{L}^i \oplus \mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R}),$$

$2 \leq k < i \leq \omega$, and

$$\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R}) \oplus \mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega).$$

In addition, clause (5) produces various sums that include summands of the form (*), such as

$$\mathcal{L}^i \oplus \mathcal{E}(\mathcal{R}, \mathcal{L}^j, \mathcal{R}) \odot \mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R})$$

TABLE I
The First Few Canonical Expressions

Generation number	\mathcal{F} -expressions	\mathcal{E} -expressions	\mathcal{A} -expressions ^a	\mathcal{G} -expressions ^a
1	\mathcal{R}	$\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{R}) = \mathcal{L}$	$\mathcal{L} \odot \dots \odot \mathcal{L} = \mathcal{L}^k$, $2 \leq k < \omega$	
2	$\mathcal{F}(\mathcal{L}) = \mathcal{L}^\omega$	$\mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R})$, $2 \leq k \leq \omega$ $\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R})$ $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega)$ $\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{L}^\omega)$	$\mathcal{L}^{k_0} \odot \mathcal{E}_1 \odot \mathcal{L}^{k_1} \odot \dots \odot$ (see text)	$\mathcal{E}(\mathcal{L}^\omega, \mathcal{R}, \mathcal{R}) \oplus$ $\mathcal{E}(\mathcal{R}, \mathcal{R}, \mathcal{L}^\omega)$ $\mathcal{L}^i \oplus \mathcal{E}(\mathcal{R}, \mathcal{L}^k, \mathcal{R})$, $2 \leq k < i \leq \omega$
.
.
.

^a This column also includes all entries in the preceding columns.

(where $j \geq 2$, $k \geq 2$, and $i > k + j$) in a profusion that makes exhaustive enumeration difficult. ■

We conclude with some open questions.

(1) A great deal of explicit information about the structure of grammatical families is contained in the canonical decomposition theorem. Can this information be used to settle some of the remaining open questions about grammatical families, such as whether the grammatical families form a lattice under inclusion? (It is not even known whether the intersection of two grammatical families is always a grammatical family.)

(2) Analogous definitions of prime and additively prime families can be made for arbitrary full semi-AFLs rather than grammatical families. Does the prime decomposition theorem still hold? What about the technical lemmas, such as 1.9? (The proof of Lemma 1.9 in the Appendix makes explicit use of the fact that the families in question are grammatical families.)

APPENDIX

The purpose of this appendix is to prove Lemma 1.9. We first prove three preliminary lemmas.

Notation. For each context-free language L , let $A(L) = \bigcap \{ \mathcal{L} \mid \mathcal{L} \text{ a grammatical family containing } L \}$.

Thus, for each infinite context-free language L , $\mathcal{S}(L) \subseteq A(L)$. No claim is made that $A(L)$ is itself a grammatical family.

A.1 LEMMA. *Let $\mathcal{S}(L)$ be an additively prime grammatical family and $L = L_1 \cup \dots \cup L_n$, where each L_i is context free. Then L is in $A(L_i)$ for some i .*

Proof. Suppose $L = \emptyset$ or $L = \{\varepsilon\}$. Then $L = L_i$ for some i and the conclusion is trivial. Suppose L is infinite and contains a word $w \neq \varepsilon$. Then w is in L_i for some i , and L is in $\mathcal{L}_{\text{fin}} = A(L_i)$. Suppose L is infinite and regular. Then some L_i is infinite, so L is in $\mathcal{R} \subseteq A(L_i)$. Finally, suppose L is not regular and the conclusion is false. Thus, L is in no $A(L_i)$. Then for each i , there is a grammatical family \mathcal{L}_i with L_i in \mathcal{L}_i and L not in \mathcal{L}_i . Since L is not regular, we may replace each trivial \mathcal{L}_i by \mathcal{R} . Hence, we may assume that each \mathcal{L}_i is nontrivial. Thus

$$\mathcal{S}(L) \subseteq \mathcal{S}(L_1) \oplus \dots \oplus \mathcal{S}(L_n) \subseteq \mathcal{L}_1 \oplus \dots \oplus \mathcal{L}_n.$$

Since $\mathcal{S}(L)$ is an additively prime grammatical family, $\mathcal{S}(L) \subseteq \mathcal{L}_i$ for some i . Then L is in \mathcal{L}_i , a contradiction. ■

A.2 LEMMA. Suppose A and B are arbitrary languages, $\#$ is a new symbol, $L \subseteq A \# B$, and L is in \mathcal{L} for some grammatical family \mathcal{L} . Then there exist $n \geq 1$ and a sequence $A_1, B_1, \dots, A_n, B_n$ of context-free languages such that:

- (1) $A_i \# B_i$ is in \mathcal{L} for each i ,
- (2) $L \subseteq \bigcup_{i=1}^n A_i \# B_i \subseteq A \# B$, and
- (3) $A \# B$ is in \mathcal{L} if A is in $\mathcal{L}(A_i)$ and B is in $\mathcal{L}(B_i)$ for some i .

Proof. Conditions (1) and (2) are slight rephrasings of parts of [7, Lemma 3.1]. The cited lemma also states that the A_i and B_i are in grammatical families $(\mathcal{L}_i)_L$ and $(\mathcal{L}_i)_R$ with the property that $A'_i \# B'_i$ is in \mathcal{L} for all A'_i in $(\mathcal{L}_i)_L$ and B'_i in $(\mathcal{L}_i)_R$ which do not have $\#$ in their alphabets. By the definition of $\mathcal{L}(A_i)$ and $\mathcal{L}(B_i)$, if A is in $\mathcal{L}(A_i)$ and B is in $\mathcal{L}(B_i)$, then A is in $(\mathcal{L}_i)_L$ and B is in $(\mathcal{L}_i)_R$, whence $A \# B$ is in \mathcal{L} . ■

A.3 LEMMA. Suppose $\mathcal{S}(A)$ and $\mathcal{S}(B)$ are additively prime grammatical families and $\#$ is a symbol not occurring in the alphabets of A and B . In addition, suppose that

$$A \# B = L \cup \bigcup_{i=1}^n A_i \# B_i,$$

where $n \geq 1$ and L, A_i , and B_i are context free for each i . Then either $A \# B$ is in $\mathcal{L}(L)$ or, for some i , A is in $\mathcal{L}(A_i)$, and B is in $\mathcal{L}(B_i)$.

Proof. Suppose there is no i , $1 \leq i \leq n$, such that A is in $\mathcal{L}(A_i)$ and B is in $\mathcal{L}(B_i)$. We shall prove that $A \# B$ is in $\mathcal{L}(L)$.

Let \mathcal{L} be an arbitrary grammatical family containing L . By Lemma A.2, there exist $m > n$ and a sequence $A_{n+1}, B_{n+1}, \dots, A_m, B_m$ of context-free languages such that $A_i \# B_i$ is in \mathcal{L} for each $i \geq n+1$, $L \subseteq \bigcup_{i=n+1}^m A_i \# B_i \subseteq A \# B$, and

- (*) $A \# B$ is in \mathcal{L} if A is in $\mathcal{L}(A_k)$ and B is in $\mathcal{L}(B_k)$ for some $k > n$.

Since $A \# B = \bigcup_{i=1}^m A_i \# B_i$, $A = \bigcup_{i=1}^m A_i$ and $B = \bigcup_{i=1}^m B_i$. Let $I = \{i \mid A \text{ not in } \mathcal{L}(A_i)\}$ and $J = \{j \mid B \text{ not in } \mathcal{L}(B_j)\}$. By Lemma A.1, $\bigcup_I A_i \neq A$, and $\bigcup_J B_j \neq B$. Select x in $A - \bigcup_I A_i$ and y in $B - \bigcup_J B_j$. Then $x \# y$ is in $A \# B$, so $x \# y$ is in $A_k \# B_k$ for some k , $1 \leq k \leq m$. Thus, x is in A_k and y is in B_k . By the choice of x and y , k is not in $I \cup J$. Hence, A is in $\mathcal{L}(A_k)$ and B is in $\mathcal{L}(B_k)$. Therefore, $n < k \leq m$. By (*), $A \# B$ is in \mathcal{L} . Since \mathcal{L} is an arbitrary grammatical family containing L , $A \# B$ is in $\mathcal{L}(L)$. ■

We are now ready to establish the desired lemma.

A.4 LEMMA (Lemma 1.9). Let $\mathcal{A}_1 \odot \mathcal{A}_2 \subseteq \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$. Then either $\mathcal{A}_1 \subseteq \mathcal{F}_1$, $\mathcal{A}_2 \subseteq \mathcal{F}_2$, or $\mathcal{A}_1 \odot \mathcal{A}_2 \subseteq \mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$.

Proof. Let $A \subseteq \Sigma^*$ and $B \subseteq \Sigma^*$ be such that $\mathcal{S}(A) = \mathcal{A}_1$ and $\mathcal{S}(B) = \mathcal{A}_2$. Let $\#$ be a new symbol. Then $A \# B$ is in $\mathcal{A}_1 \odot \mathcal{A}_2 \subseteq \mathcal{E}(\mathcal{F}_1, \mathcal{A}, \mathcal{F}_2)$. Therefore, there exist

$n \geq 1$, and for each i , $1 \leq i \leq n$, R_i in \mathcal{R} , L_i in \mathcal{A} , and \mathcal{F}_1 and \mathcal{F}_2 -substitutions τ_i and τ'_i resp., such that

$$T_i = \bigcup \{ \tau_i(w) L_i \tau'_i(w^R) / w \text{ in } R_i \}$$

and $A \neq B = T_1 \cup \dots \cup T_n$. Without loss of generality, we may assume that $L_i \neq \emptyset$, $\tau_i(a) \neq \emptyset$, and $\tau'_i(a) \neq \emptyset$ for each i and each symbol a in the alphabet of R_i . For each i , let

$$T'_i = \tau_i(R_i) L_i \tau'_i(R_i^R),$$

$$A_i = \tau_i(R_i) / (\neq \Sigma^*),$$

and

$$B_i = (\Sigma^* \#) \setminus \tau'_i(R_i^R).^6$$

Obviously, $T_i \subseteq T'_i$ is in $\mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$, since \mathcal{F}_1 and \mathcal{F}_2 are full AFL and hence closed under substitution into regular languages. Since \mathcal{F}_1 and \mathcal{F}_2 are also closed under right and left quotient by regular languages, A_i is in \mathcal{F}_1 and B_i is in \mathcal{F}_2 . Also, $A_i \subseteq A$. (For suppose x is in A_i . Then $x \neq y$ is in $\tau_i(w)$ for some y in Σ^* and w in R_i . Hence, $x \neq yz$ is in $T_i \subseteq A \neq B$ for each z in $L_i \tau'_i(w^R)$. Thus x is in A .) Similarly, $B_i \subseteq B$.

For each i , $T_i = \bigcup \{ \tau_i(w) L_i \tau'_i(w^R) \mid w \text{ in } R_i \} \subseteq \Sigma^* \# \Sigma^*$, so that either $L_i \subseteq \Sigma^* \# \Sigma^*$ or $L_i \cap \Sigma^* \# \Sigma^* = \emptyset$. Let $I = \{ i \mid L_i \subseteq \Sigma^* \# \Sigma^* \}$ and $J = \{ j \mid L_j \cap \Sigma^* \# \Sigma^* = \emptyset \}$. Suppose i is in I , i.e., $L_i \subseteq \Sigma^* \# \Sigma^*$. Let x be in $T'_i = \tau_i(R_i) L_i \tau'_i(R_i^R)$. Then $x = yz \neq z'y'$, where y is in $\tau_i(w)$ for some w in R_i , $z \neq z'$ is in L_i , and y' is in $\tau'_i(w')$, for some w' in R_i^R . Thus, for each y'' in $\tau'_i(w^R)$, $yz \neq z'y''$ is in $\tau_i(w) L_i \tau'_i(w^R) \subseteq T_i \subseteq A \neq B$. Therefore yz is in A . Similarly, $z'y'$ is in B , so x is in $A \neq B$. Hence, $T_i \subseteq T'_i \subseteq A \neq B$ for each i in I . Now suppose j is in J , i.e., $L_j \cap \Sigma^* \# \Sigma^* = \emptyset$. Then $y \neq z$ in T_j implies y in $\tau_j(R_j) / (\neq \Sigma^*) = A_j$ or z in $(\Sigma^* \#) \setminus \tau'_j(R_j^R) = B_j$. Hence, $T_j \subseteq (A_j \neq B) \cup (A \neq B_j) \subseteq A \neq B$ for each j in J . Thus, $A \neq B = \bigcup T_i = (\bigcup_I T'_i) \cup (\bigcup_J A_j \neq B) \cup (\bigcup_J A \neq B_j)$. By Lemma A.3, either

- (1) $A \neq B$ is in $\mathcal{A}(\bigcup_I T'_i)$, or
- (2) A is in $\mathcal{A}(A_j)$ for some j in J , or
- (3) B is in $\mathcal{A}(B_j)$ for some j in J .

Since each T'_i is in $\mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$, $\mathcal{A}(\bigcup_I T'_i) \subseteq \mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$. Thus (1) implies that $A \neq B$ is in $\mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2$, whence

$$\mathcal{A}_1 \odot \mathcal{A}_2 = \mathcal{S}(A) \odot \mathcal{S}(B) = \mathcal{S}(A \neq B) \subseteq \mathcal{F}_1 \odot \mathcal{A} \odot \mathcal{F}_2.$$

For each j , $\mathcal{A}(A_j) \subseteq \mathcal{F}_1$ since A_j is in \mathcal{F}_1 . Thus (2) implies that A is in \mathcal{F}_1 , and $\mathcal{A}_1 = \mathcal{S}(A) \subseteq \mathcal{F}_1$. Similarly, (3) implies that $\mathcal{A}_2 \subseteq \mathcal{F}_2$. In all cases, therefore, the conclusion of Lemma A.4 holds. ■

⁶ $R/S = \{x \mid xy \text{ is in } R \text{ for some } y \text{ in } S\}$ and $R \setminus S = \{y \mid xy \text{ is in } S \text{ for some } x \text{ in } R\}$.

REFERENCES

1. M. BLATTNER, The decidability of the equivalence of context-free grammar forms, in "20th Annual Symposium on Foundations of Computer Sciences," pp. 91-96, October 1979.
2. L. BOASSON, J. P. CRESTIN, AND M. NIVAT, Familles de langages translatables et fermées par crochet, *Acta Inform.* **2** (1973), 383-393.
3. A. CREMERS AND S. GINSBURG, Context-free grammar forms, *J. Comput. System Sci.* **11** (1975), 86-117.
4. A. CREMERS, S. GINSBURG, AND E. H. SPANIER, The structure of context-free grammatical families, *J. Comput. System Sci.* **15** (1977), 262-279.
5. S. GINSBURG, "Algebraic and Automata-Theoretic Properties of Formal Languages," North-Holland, Amsterdam, 1975.
6. S. GINSBURG, A survey of grammar forms—1977, *Acta Cybernet.* **3** (1978), 269-280.
7. S. GINSBURG, J. GOLDSTINE, AND E. H. SPANIER, A prime decomposition theorem for grammatical families, *J. Comput. System Sci.* **24** (1982), 315-361.
8. D. WOOD, "Grammar and L Forms: An Introduction," *Lecture Notes in Computer Science*, Vol. 91, Springer-Verlag, Berlin, 1980.